

This is a postprint version of the following published document:

Bega, Dario, et al. DeepCog: optimizing resource provisioning in network slicing with AI-based capacity forecasting, in *IEEE journal on selected areas in communications*, 38(2), Feb. 2020, Pp.361-376

DOI: <https://doi.org/10.1109/JSAC.2019.2959245>

© 2019 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

DeepCog: Optimizing Resource Provisioning in Network Slicing with AI-based Capacity Forecasting

Dario Bega, Marco Gramaglia, Marco Fiore, *Senior Member, IEEE*, Albert Banchs *Senior Member, IEEE*, and Xavier Costa-Perez *Senior Member, IEEE*,

Abstract—The dynamic management of network resources is both a critical and challenging task in upcoming multi-tenant mobile networks, which requires allocating capacity to individual network slices so as to accommodate future time-varying service demands. Such an anticipatory resource configuration process must be driven by suitable predictors that take into account the monetary cost associated to overprovisioning or underprovisioning of networking capacity, computational power, memory, or storage. Legacy models that aim at forecasting traffic demands fail to capture these key economic aspects of network operation. To close this gap, we present DeepCog, a deep neural network architecture inspired by advances in image processing and trained via a dedicated loss function. Unlike traditional traffic volume predictors, DeepCog returns a cost-aware *capacity forecast*, which can be directly used by operators to take short- and long-term reallocation decisions that maximize their revenues. Extensive performance evaluations with real-world measurement data collected in a metropolitan-scale operational mobile network demonstrate the effectiveness of our proposed solution, which can reduce resource management costs by over 50% in practical case studies.

Index Terms—Mobile networks, Network Slicing, 5G networks, Artificial Intelligence, Deep Learning.

I. INTRODUCTION

Network slicing is an emerging paradigm that is expected to characterize 5G and beyond-5G mobile systems. Network slicing allows operators to tailor Virtual Network Functions (VNFs) to the precise requirements of individual mobile services [1], and will be key in enabling the unprecedented heterogeneity of future mobile applications [2]. Managing sliced networks will represent a major challenge for operators, since this new paradigm represents a shift from the rather limited reconfiguration possibilities offered by current Operations and Business Support System (OSS/BSS) to a complex, software-defined control layer that must dynamically organize thousands of slices belonging to hundreds of tenants on the same infrastructure [3].

D. Bega and A. Banchs are with the IMDEA Networks Institute, 28918 Leganés, Spain, and also with the Telematic Engineering Department, University Carlos III of Madrid, 28911 Leganés, Spain (e-mail: dario.bega@imdea.org; banchs@it.uc3m.es).

M. Gramaglia is with the Telematic Engineering Department, University Carlos III of Madrid, 28911 Leganés, Spain (e-mail: mgramagl@it.uc3m.es).

M. Fiore is with CNR, 10129 Torino, Italy. (e-mail: marco.fiore@ieiit.cnr.it)

X. Costa-Pérez are with NEC Laboratories Europe, 69115 Heidelberg, Germany (e-mail: xavier.costa@neclab.eu).

Manuscript received June 22, 2019; revised September 3, 2019; accepted November 6, 2019. Date of publication xxxx xx, xxxx; date of current version xxxx xx, xxxx. (Corresponding author: Dario Bega.)

A. Network management and capacity forecast.

To rise to the challenge, network operators must introduce substantial automation in the presently human-driven management and orchestration (MANO) processes, ultimately realizing the 5G principle of *cognitive network management* [4]. Achieving this objective requires advances in two complementary technologies: (i) technical solutions that enable end-to-end Network Function Virtualization (NFV), providing the flexibility necessary for resource reallocation; and, (ii) a network intelligence that automatically identifies and anticipates demand patterns from streaming network measurement data, and then takes decisions on how to configure VNFs and allocate resources so as to maximize the system efficiency.

From a technical standpoint, solutions that implement NFV at different network levels are well established, and start to be tested and deployed. Examples include current MANO platforms architectures like ETSI NFV [5], and implementations of MANO controllers such as OSM [6] or ONAP [7] that support reconfiguring resources to VNFs on the fly. By contrast, the integration of intelligence in cognitive mobile networks is still at an early stage. Nowadays, resource assignment to VNFs is a reactive process, mostly based on hysteresis thresholding and aimed at self-healing or fault tolerance. There is a need for proactive, data-driven, automated solutions that enable cost-efficient network resource utilization, by anticipating future needs for capacity and timely reallocating resources just where and when they are required.

The aim of our work is precisely to contribute to the definition of a network intelligence that is adapted to a network slicing environment. More specifically, the focus of this paper is on the design of data analytics that enable the anticipatory allocation of resources in cognitive mobile networks. To this end, we investigate a machine learning solution that runs on traffic measurement data and provides operators with information about the capacity needed to accommodate future demands at each network slice – a critical knowledge for data-driven resource orchestration. We take a pragmatic approach, and duly account for the economic costs associated to the operation above.

Indeed, resource orchestration decisions have a direct monetary impact for the network operator in terms of operating expenses, which can be divided into two macroscopic categories of cost.

- *Overprovisioning* – when providing excess capacity with

respect to the actual resource demand, the operator incurs a cost due to the fact that it is reserving more resources than those needed to a network entity (*e.g.*, a network slice, a network function, or a virtual machine). As resources are typically isolated across slices, this seizes the excess resources from other network entities that may have possibly used them. At a global system level, continued overprovisioning implies that the operator will have to deploy more resources than those required to accommodate the user demand, limiting the advantage of a virtualized infrastructure and of cognitive networking solutions in general.

- *SLA violation* – if insufficient resources are allocated to a network entity, users will suffer low Quality of Service (QoS), or even discontinued service. This has an indirect price for the operator, in terms of customer dissatisfaction and increased churning rates, which is not simple to quantify. However, in emerging contexts such as those promoted by network slicing, underprovisioning also entails a different, more direct and quantifiable economic penalty for the operator. Under slicing, operators will sign Service Level Agreements (SLAs) with the mobile service providers, which need to be strictly enforced. Underprovisioning means violating such SLAs, which results in substantial monetary fees for the network operator.

Clearly, the cost is not the same in the two cases, and it may also vary depending on the specific settings, including the nature of the concerned resources, the technologies deployed in the network infrastructure, or the market strategies of the operator. In all cases, we posit that, once suitably modeled, such costs shall be at the core of the orchestrating decisions.

Legacy techniques for the prediction of mobile network traffic, such as the one reviewed in Section II, fall short in this respect. Such models aim at perfectly matching the temporal behavior of traffic, independently of whether the anticipated demand is above or below the target, and are thus agnostic of the aforementioned costs. As a result, they return forecasts as that depicted in Fig. 1a, which refers to a real-world case study of YouTube video streaming traffic at a core network datacenter. Note that no distinction is made between positive and negative errors, which leads to substantial SLA violations covering roughly half of the observation time. The operator may then attempt to apply overprovisioning to the output provided by such a traffic predictor. Unfortunately, legacy forecast models do not offer any insight on how large the excess resource allocated on top of the forecast demand should be. As we will demonstrate later in the paper, this makes such a strategy highly inefficient.

B. Paper contributions and data-driven setup

In this paper, we present DeepCog, a new mobile traffic data analytics tool that is explicitly tailored to solve the capacity forecast problem exposed above. The design of DeepCog yields multiple novelties, summarized as follows:

- It hinges on a deep learning architecture inspired by recent advances in image and video processing, which

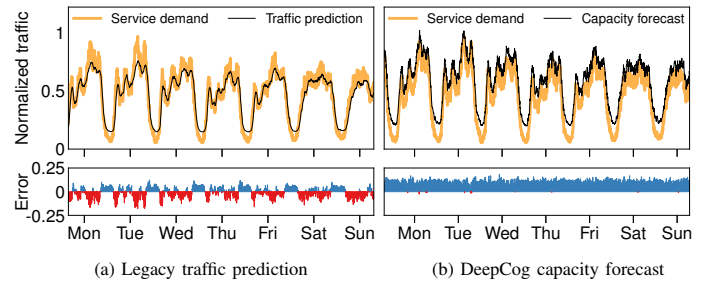


Fig. 1. Top: actual and predicted weekly demands for YouTube at a datacenter controlling 470 4G eNodeBs. Bottom: levels of overprovisioning (blue) and capacity violations (red) over time. (a) Output of a recent deep learning predictor of mobile traffic [8]. (b) Output of DeepCog, tailored to anticipatory network resource allocation. The best view of this figure is in colors.

exploits space- and time-independent correlations typical of mobile traffic and computes outputs at a datacenter level;

- It leverages a customized loss function that targets capacity forecast rather than plain mobile traffic prediction, letting the operator tune the balance between overprovisioning and demand violations;
- It provides long-term forecasts over configurable prediction horizons, operating on a per-service basis in accordance with network slicing requirements.

Overall, these design principles jointly solve the problem of capacity forecast in network slicing. This is illustrated by Fig. 1b, which shows an example of the required capacity forecast by DeepCog in a real-world case study. We remark that DeepCog is one of the very first examples of rigorous integration of machine learning into a cognitive network management process, and marks a difference from the common practice of embedding vanilla deep learning structures into network operation [9].

Extensive performance evaluations with substantial measurement data collected in an operational metropolitan-scale mobile network demonstrate the superiority of our approach over a wide range of benchmarks based on traditional and state-of-the-art mobile traffic predictors. Even though our performance evaluation focuses on cellular network environments, it is worth noting that the data-driven approach adopted is highly flexible and can be applied to forecast capacity in different scenarios, including, *e.g.*, other types of access networks.

The document is organized as follows. We first provide a review of related works, and highlight the novelty of our proposed method, in Section II. We then outline the overall framework of DeepCog in Section III, and detail the design of its most critical component, *i.e.*, the loss function, in Section IV. The quality of the solution is then assessed in realistic scenarios in Section V. Finally, we draw conclusions in Section VI.

II. RELATED WORK

Applications to networking problems of machine learning in general, and of deep learning in particular, are starting to become popular. Artificial intelligence can indeed be applied to solve many different problems that emerge in computer

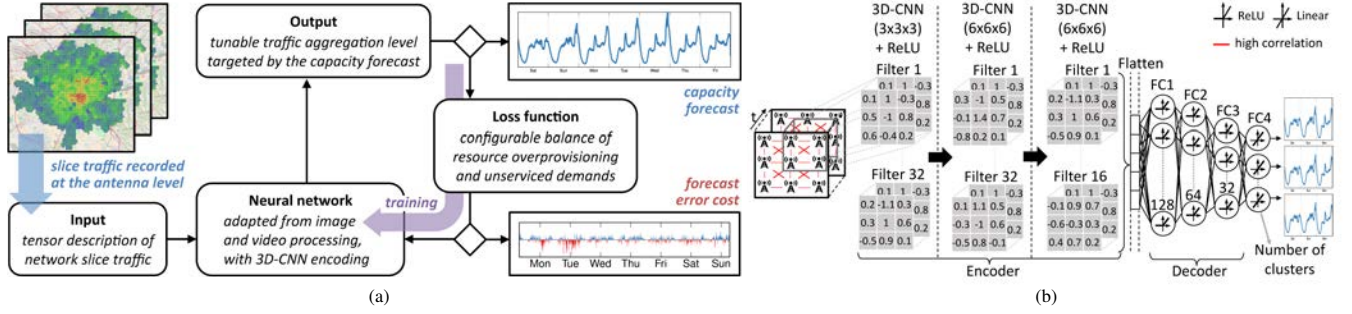


Fig. 2. The overall DeepCog architecture. (a) Outline and interaction of the DeepCog components. (b) DeepCog neural network encoder-decoder structure.

networks, as highlighted in recent comprehensive surveys on the topic [9], [10].

In the context of network management, emerging paradigms like slicing increase substantially the complexity of orchestrating network functions and resources, at all levels. For instance, intelligence is needed for the admission control of new slices: as resources are limited and slicing entails their strong isolation, this is critical to ensure that the system operates efficiently. With potentially hundreds of slices allocated simultaneously, and a need to anticipate highly profitable future requests, the decision space for admission control becomes so large that traditional approaches become impractical. Solutions based on deep learning architectures represent here a viable approach [11]. Similar considerations apply to other aspects of sliced network management, *e.g.*, the allocation of computational resources to slices at the radio access, based on transmission (*e.g.*, modulation and coding scheme, channel load) and environmental (*e.g.*, signal quality, hardware technology) conditions [12], or the anticipatory reservation of Physical Resource Blocks (PRBs) to user traffic to be served in target network slices [13].

Our specific problem relates to the orchestration of generic resources (*e.g.*, CPU time, memory, storage, spectrum) to slices at different network entities, which is tightly linked to mobile traffic prediction. The literature on forecasting network traffic is in fact vast [9], [14]. Solutions to anticipate future offered loads in mobile networks have employed a variety of tools, from autoregressive models [15]–[17] to information theoretical tools [18], passing by Markovian models [19] and deep learning [8], [11], [13], [20], [21].

However, we identify the following major limitations of current predictors when it comes to supporting resource orchestration in mobile networks.

First, predictors of mobile traffic invariably focus on providing forecasts of the future demands that minimize some absolute error [9], [14]. This approach leads to predicted time series that deviate as little as possible from the actual traffic time series, as exemplified in Fig. 1a for a real-world case study. While reasonable for many applications, such an output is not appropriate for network resource orchestration. As explained in Section I, the operator aims at provisioning sufficient capacity to accommodate the offered load at *all* times, since failing to do so implies high costs in terms of high subscribers' churn rates, as well as significant fees for

violating SLAs signed with tenants. Yet, if an operator decided to allocate resources based on a legacy prediction like that in Fig. 1a, it would incur into capacity violations most of the time (as illustrated in the bottom subplot).

Second, with the adoption of network slicing, forecasts must occur at the slice level, *i.e.*, for specific mobile services in isolation. However, most traffic predictors, including recent ones, are evaluated with demands aggregated over all services [8], [20], [21]. This is an easier problem, since aggregate traffic yields smoother and more regular dynamics, hence previous solutions may not handle well the bursty, diversified traffic exhibited by each service. The only attempts at anticipating the demands generated by specific mobile services have been made by using multiple-input single-output (MISO) autoregressive models [22], and hybrid prediction methods that incorporate α -stable models and sparsity with dictionary learning [18].

Third, existing machine learning predictors for mobile traffic typically operate at base station level [8], [21]. However, NFV operations mainly occur at datacenters controlling tens (*e.g.*, at the mobile edge) to thousands (*e.g.*, in the network core) of base stations. Here, prediction should be more efficient when performed on the aggregate traffic at each datacenter, where orchestration decisions are taken, rather than combining independent forecasts from each base station.

Our proposed solution, DeepCog, addresses all of the open problems above, by implementing a first-of-its-kind predictor that anticipates the minimum provisioned capacity needed to cut down SLA violations. This closes the present gap between traffic prediction and practical orchestration, as it provides the operator with an explicit capacity forecast that mitigates underprovisioning in Fig. 1b while minimizing unnecessary resource reservation. We remark that early versions of the DeepCog framework were presented in [23] and [24]. Those preliminary variants of our solution could achieve short-term capacity forecasting over the next time step, whereas the complete version presented in this paper supports long-term capacity prediction over a configurable number of future time steps.

III. A DEEP LEARNING FRAMEWORK FOR RESOURCE ORCHESTRATION

The design of DeepCog is outlined in Fig. 2a. Its organization is that typical of deep learning systems, and it stems from

(i) properly formatted *input* data used to build the forecast, which, in our case, represents the current and past traffic associated to a specific network slice as a tensor. Such input is fed to (ii) a *deep neural network* architecture that extrapolates and processes input features to provide (iii) an *output* value: the capacity forecast. During the training phase, the output is used to evaluate (iv) a *loss function* that quantifies the error with respect to the ground truth, and, in DeepCog, accounts for the costs associated to resource overprovisioning and service request denial.

In our network model, we consider that time is divided in slots, which we denote by t . Let $\delta_s^i(t)$ be the traffic associated with slice s that is observed at base station $i \in \mathcal{N}$ and time t . A *snapshot* of the demand of slice $s \in \mathcal{S}$ at time t is given by a set $\delta_s(t) = \{\delta_s^1(t), \dots, \delta_s^N(t)\}$, and provides a global view of the traffic for that slice at time t across the whole network. We let \mathcal{N} denote the set of N base stations in the network, and \mathcal{M} the set of $M < N$ datacenters. For each slice $s \in \mathcal{S}$, base stations are associated to datacenters via a surjective mapping $f_s : \mathcal{N} \rightarrow \mathcal{M}$, such that a datacenter $j \in \mathcal{M}$ serves the aggregated load of slice s for all of the associated bases stations¹. With this mapping, the traffic for slice s processed by datacenter j at time t is given by $d_s^j(t) = \sum_{i|f_s(i)=j} \delta_s^i(t)$. Then, the set of demands across all datacenters is given by $\mathbf{d}_s(t) = \{d_s^1(t), \dots, d_s^M(t)\}$.

Let us denote the allocated capacity for slice s at datacenter j and time t as $c_s^j(t)$, and the set of capacities at all $j \in \mathcal{M}$ as $\mathbf{c}_s(t) = \{c_s^1(t), \dots, c_s^M(t)\}$. Then, the capacity forecast problem is that of computing a *constant* capacity $\bar{\mathbf{c}}_s(t, T_h) = \{\bar{c}_s^1(t, T_h), \dots, \bar{c}_s^M(t, T_h)\}$ that is allocated in the network datacenters over a time horizon T_h , i.e., through an interval between the present time t and a future time $t + T_h$. In practice, this models the typical situation where the resource reconfiguration frequency is limited (e.g., by the NFV technology), and the operator must decide in advance the amount of resources that will stay assigned to a slice until the next reallocation takes place. The time horizon T_h thus corresponds to the reconfiguration period, and the allocated capacity is such that $c_s^j(t) = \bar{c}_s^j(t, T_h) \forall j \in \mathcal{M}, t \in [t, t + T_h]$.

The forecast builds on knowledge of the previous T_p traffic snapshots $\delta_s(t-1), \dots, \delta_s(t-T_p)$. The quality of the capacity forecast $\bar{\mathbf{c}}_s(t, T_h)$ is measured by means of a suitable loss function $\ell(\bar{\mathbf{c}}_s(t, T_h), \mathbf{d}_s(t), \dots, \mathbf{d}_s(t + T_h))$. This function $\ell(\cdot)$ determines the compound cost of overprovisioning and underprovisioning network resources at the target datacenters, as produced by allocating a constant capacity $\bar{\mathbf{c}}_s(t, T_h)$ when the actual time-varying demand is in fact $\mathbf{d}_s(t), \dots, \mathbf{d}_s(t + T_h)$.

Below, we present each of the components of the framework, and discuss its mapping to the elements of a 5G network architecture running cognitive resource management.

A. The Neural Network

DeepCog leverages a deep neural network structure composed of suitably designed encoding and decoding phases, performing a capacity forecasting prediction over a given time horizon. The structure is general enough that it can be trained to solve the capacity forecast problem for (i) network slices dedicated to different services with significantly diverse demand patterns, (ii) any datacenter configuration, and (iii) any time horizon T_h . The hyperparameters of the neural network have been tuned through extensive simulation and testing.

The design of the neural network structure in DeepCog is inspired by recent breakthroughs [25] in deep learning for image and video processing. As summarized in Fig. 2b, the network is composed of an encoder that receives an input representing the mobile traffic data $\delta_s(t-1), \dots, \delta_s(t-T_p)$ and maps important spatial and temporal patterns in such data onto a low-dimensional representation. Intuitively, the encoder extracts the relevant features from the input traffic tensors $\delta_s(t-1), \dots, \delta_s(t-T_p)$ and the decoder leverages such features to generate a capacity forecast that is tailored to a given combination of slice, prediction time horizon, and datacenter class; e.g., datacenters deployed close to the radio access will show different features from those co-located with the Internet gateways.

The result of the encoder undergoes a flattening process that converts the 3D (space and time) tensor data into a unidimensional vector format. This is the input format required by the fully connected layers that form the decoder, which then generates the final capacity forecast $\bar{\mathbf{c}}_s(t, T_h)$ at the target set of datacenters \mathcal{M} . Below, we detail the encoder and decoder implementations, and discuss the training procedure.

1) *The Encoder*: The encoder is composed by a stack of three three-dimensional Convolutional Neural Network (3D-CNN) layers [26]. Generic Convolutional Neural Networks (CNNs) are a specialized kind of deep learning structure that can infer local patterns in the feature space of a matrix input. In particular, two-dimensional CNNs (2D-CNNs) have been extensively utilized in image processing, where they can complete complex tasks on pixel matrices such as face recognition or image quality assessment [27]. 3D-CNNs extend 2D-CNNs to the case where the features to be learned are spatiotemporal in nature, which adds the time dimension to the problem and transforms the input into a 3D-tensor. Since mobile network traffic exhibits correlated patterns in both space and time, our encoder employs 3D-CNN layers².

Formally, the 3D-CNN layers receive a tensor input $\mathcal{T}(\delta_s(t-1)), \dots, \mathcal{T}(\delta_s(t-T_p))$, where $\mathcal{T}(\cdot)$ is a transformation of the argument snapshot into a matrix. This input is processed by three subsequent 3D-CNN layers. Each neuron of these layers runs a filter $\mathcal{H}(\sum_{\tau} \mathbf{I}(\tau) * \mathbf{K}(\tau) + \mathbf{b})$ where $\mathbf{I}(\tau)$ is the input matrix passed to the neuron (e.g., $\mathbf{I}(\tau) =$

¹We remark that DeepCog works for any arbitrary mapping, including, e.g., flows from a slice in the same base station being split across datacenters, or associations among base stations and datacenters varying over time. As a matter of fact, DeepCog's learning process is based exclusively on the traffic load at each individual datacenter and is thus independent of the actual sources generating such traffic.

²We have employed CNNs instead of Recurrent Neural Networks (RNNs) (typically used for forecasting application) because the mobile load at a time instant t mainly depends on previous T_p instants and not on all the past values (as confirmed by our analysis). For this reason, CNNs provide us enough temporal memory while being cheaper to train in terms of computational cost compared with RNNs.

$\mathcal{T}(\delta_s(\tau))$ at the very first layer, for slice s and generic time τ , $*$ denotes the 3D convolution operator, $\mathbf{K}(t)$ is the kernel of filters, $\mathcal{H}(\cdot)$ is a non-linear activation function, and \mathbf{b} is a bias vector. We use two different kernel configurations $\mathbf{K}(\tau)$, as shown in Fig. 2b: a $3 \times 3 \times 3$ kernel for the first 3D-CNN layer, and a $6 \times 6 \times 6$ kernel for the second and third layers. These settings allow limiting the *receptive field*, i.e., the portion of input analyzed by each neuron, to small regions: in presence of strong local correlation of the input data, this approach is known to yield good performance with fairly limited training, in particular compared to Recurrent Neural Networks. As for the choice of the activation function, many different options have been proposed in the literature, spanning from linear functions to tanh, sigmoid or Rectified Linear Unit (ReLU). Among these, we select ReLU, and set $\mathcal{H}(\mathbf{x}) = \max(0, \mathbf{x})$, which provides advantages in terms of discriminating performance and faster learning [28]. Finally, \mathbf{b} is randomly set at the beginning of each training phase.

The second and third 3D-CNN layers are interleaved with Dropout layers: such layers regularize the neural network and reduce overfitting [28] by randomly setting to zero a number of output features from the preceding layer during the training phase. The *dropout rate* defines the probability with which output features undergo this effect. During training, we employ two Dropout layers with dropout rate equal to 0.3.

2) *The decoder*: The decoder uses Multi-Layer Perceptrons (MLPs) [29], a kind of fully-connected neural layers, where every neuron of one layer is connected to every neuron of the next layer. This provides the ability to solve complex function approximation problems. In particular, MLPs are able to learn global patterns in their input feature space [30], allowing the neural network structure to forecast the targeted load value leveraging the local features extracted by the encoder. In our structure, each layer performs an operation $\mathcal{H}'(\mathbf{x} \cdot \mathbf{W} + \mathbf{b})$, where \mathbf{x} is the MLP layer input vector, \mathbf{W} a weight matrix related to the neurons of each layer, and \mathbf{b} the bias vector. \mathbf{W} plays a similar role to $\mathbf{K}(t)$ in the encoder part: its values drive the prediction through the layers of the decoding part.

As for the activation functions \mathcal{H} , we employ ReLU for all MLP layers except for the last one, where a linear activation function is used since the desired output takes real values. The last linear layer can be configured to produce multiple predictions in parallel, each matching the aggregate capacity required by a subset of base stations, thus allowing to forecast the needed capacity for different datacenters comprising a subset of base stations. Ultimately, this organization makes the DeepCog neural network capable of predicting per-slice capacity requirements at datacenter level, in a way that can adapt to any configuration of \mathcal{M} and to any time horizon T_h .

3) *The training procedure*: We leverage the popular Adam optimizer, which is a Stochastic Gradient Descent (SGD) method that provides faster convergence compared to other techniques [31]. SGD trains the neural network model, evaluating at each iteration the loss function $\ell(\cdot)$ between the forecast and the ground truth, and tuning the model parameters in order to minimize $\ell(\cdot)$. For the configuration of the Adam optimizer, we use the default configuration with a learning rate of 5×10^{-4} .

An important element that concerns the training of the DeepCog architecture is that the encoder and the decoder described in Section III-A1 and Section III-A2 have independent roles. Therefore, while the decoder heavily depends on the forecast specifications, the encoder does not, and is agnostic to the final usage of the extracted features. This fact allows adopting a *transfer learning* approach during training: instead of treating the two blocks as a whole (and performing the training over the full system for all the possible slices, datacenter classes and horizons), we can train them separately. Specifically, an horizon-independent encoder can be trained on past traffic tensors at maximum time granularity, and then reused in combination with dedicated decoders tailored to each T_h value. Beside reducing the training time, this strategy reduces the need for neural-network-wide training to different settings of slice and datacenter only.

B. Arrangement of input data

The input is composed by measurement data generated in a specific network slice, and recorded by dedicated probes deployed within the network infrastructure. Depending on the type and location of the probe, the nature of the measurement data may vary, describing the demands in terms of, e.g., signal quality, occupied resource blocks, bytes of traffic, or computational load on VNFs. DeepCog leverages a set of transformations to map any type of slice traffic measurements into a tensor format that can be processed by the learning algorithm.

The 3D-CNN layer adopted as the first stage of the decoder requires a multidimensional tensor input. We thus need to define the transformation $\mathcal{T}(\cdot)$ of each traffic snapshot into a matrix. Note that 3D-CNN layers best perform in presence of a tensor input that features a high level of local correlation, so that neurons operate on similar values. In image processing, where close-by pixels typically have high correlation, this is easily solved by treating the pixel grid as a matrix. In line with this strategy, the current common practice in mobile network traffic prediction is to leverage the geographical locations of the base stations, and assign them to the matrix elements so that their spatial proximity is preserved as much as possible [8], [9]. However, this approach does not consider that correlations in mobile service demands at a base station level do not depend on space, rather on land use [32]: base stations exhibiting strongly correlated network slice traffic may be far apart, e.g., covering the different train stations within a same large city. Thus, we aim at creating a tensor input whose neighboring elements correspond to base stations with strongly correlated mobile service demands. To this end, we construct the mapping of base stations into a matrix structure as follows.

- For each base station i , we define its historical time series of total traffic as $\tau^i = \{\delta^i(1), \dots, \delta^i(t-1)\}$, where $\delta^i(t) = \sum_s \delta_s^i(t)$. Then, for each pair i and j , we determine the similarity of their recorded demands by computing $\text{SBD}^{ij} = f_{\text{SBD}}(\tau^i, \tau^j)$, where $f_{\text{SBD}}(\cdot)$ is the shape-based distance, a state-of-the-art similarity measure for time series [33]. All pairwise distances are then stored in a distance matrix $\mathbf{D} = (\text{SBD}^{ij}) \in \mathbb{R}^{M \times M}$.

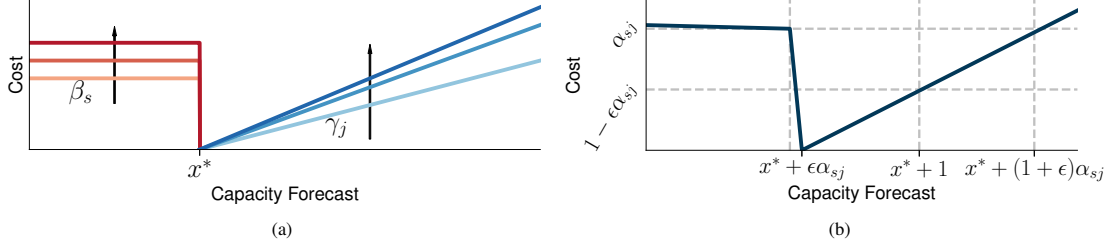


Fig. 3. Graphical representations of the α -OMC loss function. (a) Ideal model. (b) Actual implementation in (1).

- We compute virtual bidimensional coordinates \mathbf{p}_i for each base station i so that the values in the distance matrix \mathbf{D} are respected as much as possible. Formally, this maps to an optimization problem whose objective is $\min_{x_1, \dots, x_M} \sum_{i < j} (\|\mathbf{p}_i - \mathbf{p}_j\| - \text{SBD}^{ij})^2$, efficiently solved via Multi-Dimensional Scaling (MDS) [34].
- We match each point \mathbf{p}_i to an element e of the input matrix \mathbf{I} , again minimizing the total displacement. To this end, we: (i) quantize the virtual surface encompassing all points \mathbf{p}_i so that it results into a regular grid of N cells; (ii) assume that each cell is an element of the input matrix; (iii) compute the cost k_{ie} of assigning a point \mathbf{p}_i to element e as the Euclidean distance between the point and the cell corresponding to e . We then formalize an assignment problem with objective $\min_a \sum_{i \in \mathcal{N}} \sum_{e \in \mathbf{I}} k_{ie} x_{ie}$, where $x_{ie} \in [0, 1]$ is a decision variable that takes value 1 if point \mathbf{p}_i is assigned to element e , and must fulfill $\sum_{i \in \mathcal{N}} x_{ie} = 1$ and $\sum_{e \in \mathbf{I}} x_{ie} = 1$. The problem is solved in polynomial time by the Hungarian algorithm [35].

The solution of the assignment problem is the transformation $\mathcal{T}(\cdot)$ of the original base stations into elements of the matrix \mathbf{I} . The mapping function $\mathcal{T}(\cdot)$ allows translating a traffic snapshot $\delta_s(t)$ into matricial form. Applying this to snapshots at different times, $\delta_s(t-1), \dots, \delta_s(t-T)$, we can thus build the tensor required by the entry encoder layer in Fig. 2b.

C. The Output function

DeepCog is designed for flexibility, and can be used for different orchestration scenarios. This is achieved thanks to an adaptable last layer of the deep neural network, and a configurable loss function. In general, the learning algorithm returns a forecast of the capacity required to accommodate the future demands for services associated to a specific network slice. This generic definition of output can then be applied to different orchestration use cases that may differ in the traffic aggregation level at which the resource configuration takes place, and/or in the frequency at which resource reallocation can be realized.

For instance, the anticipatory assignment of baseband processing units to network slices in a Cloud Radio Access Network (C-RAN) datacenter requires a prediction of the capacity needed to accommodate the traffic of a few tens of base stations; instead, reserving memory resources for a specific network slice at a core network datacenter implies forecasting capacity for the data sessions of subscribers associated to hundreds of base stations. The output format of DeepCog can

accommodate any datacenter layout, by tailoring the last linear layer of the neural networks to the specific requirements of the layout (as discussed in Section III-A2).

Also, as discussed previously, the time horizon over which the forecast is performed is another relevant system parameter, which depends on NFV technology limitations and current trends in *commoditization* of softwarized mobile network. When the technology limitations do not allow frequent reconfiguration opportunities, resources need to be allocated over long periods, e.g., of tens of minutes or even hours. In this case, forecasting over long-term horizons provides the operator with information on the constant capacity to be allocated during long intervals. To realize this, DeepCog operates on configurable time horizons, thanks to the flexible loss function that we will discuss next.

IV. α -OMC

One of the key components of the system proposed in the previous section is the *loss function*, denoted by $\ell(\cdot)$. This function determines the penalty incurred when making a prediction error. In this paper, we propose a novel loss function that is tailored to the specific requirements of the capacity forecast problem. Our design of $\ell(\cdot)$ accounts for the costs resulting from (i) forecasting a lower value than the actual offered load, which leads to an *SLA violation* due to the provisioning of insufficient resources, (ii) predicting a higher value than the actual one, which leads to *overprovisioning*, allocating more resources than those needed to meet the demand. In order to ensure that we drive the system towards an optimal trade-off between overprovisioning and SLA violations, over a generic time horizon T_h , $\ell(\cdot)$ must account for the penalty inflicted in each case. In what follows, we describe the design of α -OMC (Operator Monetary Cost), a loss functions that provides DeepCog with the capability of optimizing the overall running costs of the system.

A. Loss function design

In DeepCog, the loss function steers the behavior of the neural network by adjusting the weights of the neurons according to the error between the estimated value and the real one. To achieve the objective of minimizing the overall cost, a custom loss function for the capacity forecasting problem is composed by a term $f(x, x^*)$ that deals with the resource overprovisioning penalty, and a term $g(x, x^*)$ that models the cost of resource violations. The variable x represents the

allocated resources at a given time interval, while x^* is the real demanded load for the same period. So the overall cost is due by the discrepancy between x and x^* in any time horizon.

The shape of overall cost function $f(x, x^*) + g(x, x^*)$ is depicted in Fig. 3a. A perfect algorithm (*i.e.*, an *oracle*) always keeps the system in the optimal operation point $x = x^*$ where no penalty is introduced, *i.e.*, $f(x^*, x^*) = g(x^*, x^*) = 0$. Of course, errors are inherent to predictions, and it is very unlikely that the forecast perfectly matches the real demand: hence, a penalty value is back-propagated depending on whether x is above or below the target operation point x^* .

1) $g(x, x^*)$, a *reactive approach to SLA violations*: When the orchestrated resources are less than those needed in reality (*i.e.*, $x < x^*$) the network operator pays a monetary compensation to the tenant. We assume an SLA that enforces a proportional compensation depending on the number of time intervals in which an operator fails to meet the requirements set by a tenant due to insufficient capacity allocated to the slice. Thus, SLA violations determine a fixed cost for the operator at every time interval where the tenant demand is not satisfied. Accordingly, we let the system learn that the operation point x^* is actually higher than the estimated one by applying a penalty β_s as soon as the estimation falls below the real value. The parameter β_s can be customized to the needs of the slice $s \in \mathcal{S}$: higher values may be used for cases where reliability is paramount like, *e.g.*, in URLLC network slices; instead, lower values can be applied for slices where KPI commitments are provided over longer time intervals. Note that higher β_s values are likely to bring the system toward $x > x^*$, incurring hence in higher deployment costs, as discussed next.

2) $f(x, x^*)$, a *monotonically increasing cost for resource overprovisioning*: While SLA violations depend on the agreements between the tenants and the operator, the overprovisioning cost solely depends on the network operator, and more specifically on the deployment costs associated with excess allocated capacity. We assume that such a cost grows with the amount of unused capacity at each time interval, and model it as a positive monotonic function that is only applied when $x > x^*$: the higher the resource provisioning error, the more (unnecessarily) expensive is the deployment. The exact expression of $f(x, x^*)$ may vary, and one could consider, *e.g.*, linear, super-linear, or exponential shapes. For DeepCog, we design α -OMC to use a linear function, as shown in Fig. 3a. The linear scaling factor γ_j is configurable by the operator, and represents the monetary cost of the excess resource allocation. The cost depends on the specific datacenter $j \in \mathcal{M}$ at which the capacity forecasting takes place: for instance, spectrum resources at the edge are typically scarcer and more expensive to deploy than computational resources in a network core datacenter. In case of expensive resources (characterized by a large γ_j), a positive forecasting error will have a higher impact, favoring a capacity forecast with a lower level of overprovisioning.

3) *Balancing the two cost contributions.*: Overall, the amount of resources that a network operator is willing to allocate depends on the cost that it has to pay when failing to meet the demands for a given slice (given by β_s) and the

cost associated with adding extra resources at a specific data-center (given by γ_j). These two parameters, β_s and γ_j , push the capacity allocations towards opposite directions, namely overdimensioning and underdimensioning, respectively. Rather than their absolute values, what really matter for the resulting allocation is the ratio between the two parameters, which determines the trade-off between overdimensioning and underdimensioning. Accordingly, in the following we express the custom loss as a function of a single parameter $\alpha_{sj} \doteq \frac{\beta_s}{\gamma_j}$. We remark that α_{sj} indicates the monetary costs of SLA violations with respect to the overprovisioning: failing to meet the slice requirements once costs as much as allocating α_{sj} units of excess capacity. Thus, a higher α_{sj} implies higher SLA violation costs relative to the deployment (*i.e.*, overprovisioning) cost. A mobile network operator can easily set this parameter based on its deployment costs, SLA fees, and market strategies.

Another important remark is that the SGD method used to train the neural network does not work with constant or step functions, and requires that the loss function be differentiable in all its domain. We solve this problem by introducing minimum slopes of very small intensity ϵ for $x < x^*$ and at $x = x^*$. We name the resulting loss function Operator Monetary Cost, which has a single configurable parameter α_{sj} . The final expression of α -OMC is

$$\alpha\text{-OMC}(x, x^*) = \begin{cases} \alpha_{sj} - \epsilon(x - x^*) & \text{if } x \leq x^* \\ \alpha_{sj} - \frac{1}{\epsilon}(x - x^*) & \text{if } x^* < x \leq x^* + \epsilon\alpha_{sj} \\ x - x^* - \epsilon\alpha_{sj} & \text{if } x > x^* + \epsilon\alpha_{sj}. \end{cases} \quad (1)$$

Fig. 3b provides a sample illustration of (1) above.

The final loss function $\ell(\cdot)$ then measures the quality of the forecast over the time horizon T_h , by applying the α -OMC expression over multiple time intervals as follows:

$$\ell(\bar{\mathbf{c}}_s(t, T_h), \mathbf{d}_s(t), \dots, \mathbf{d}_s(t + T_h)) = \sum_{j \in \mathcal{M}} \sum_{\tau=0}^{T_h} \alpha\text{-OMC}(\bar{c}_s^j(t, T_h), d_s^j(t + \tau)). \quad (2)$$

For the sake of readability and without loss of generality, in the remainder of the paper we will employ a constant $\beta_s = \beta$ across slices and a constant $\gamma_j = \gamma$ across datacenter deployment, leading to $\alpha_{sj} = \alpha$ for all $s \in \mathcal{S}, j \in \mathcal{M}$.

B. Correctness and convergence

We now analyze the proposed loss function in terms of (i) correctness, *i.e.*, its capability of achieving a performance that is close to the optimal, and (ii) convergence, *i.e.*, the time it requires to learn such a correct strategy.

In Fig. 4, we run DeepCog in the representative network resource management case studies that are later detailed in Section V, where a slice is dedicated to one particular mobile service and runs in a specific class of network datacenter. For each case study, DeepCog forecasts a given level of capacity to be allocated which leads to an associated monetary cost. In order to investigate the correctness of the solution, we vary

the provisioned capacity by adding to or subtracting a fixed offset from the capacity indicated by DeepCog.

The curves of Fig. 4 illustrate the variation of the monetary cost (in the y axis) as the offset is shifted (in the x axis), where increasingly positive (respectively, negative) values on the x axis correspond to a higher (respectively, lower) level of capacity provisioning with respect that suggested by our solution. The results prove that DeepCog always identifies the capacity allocation that minimizes the monetary cost for the operator under the inherently inaccurate prediction, as both a higher and a lower level of overprovisioning leads to a greater cost. This holds under any combination³ of target mobile service, datacenter class, and system settings α or T_h , which demonstrates the high consistency of our solution in balancing costs caused by SLA violations and overprovisioning.

We next assess the convergence properties of the loss function that drives DeepCog, by observing its behavior over time. Specifically, we measure the normalized cost of the solution identified by our learning algorithm, and compare it against that returned by the same neural network trained with legacy loss functions.

Fig. 5 shows how the average normalized cost of network operation varies during the training phase for different α , services and datacenter classes. While the α -OMC loss function minimizes the monetary cost of the operator in less than 20 epochs, both MAE and MSE converge to a fixed fee that grows as α increases. This confirms that classical loss functions are not effective when dealing with capacity forecasting, resulting in high penalties for operators. The results are consistent across all of the different configuration scenarios we tested.

V. PERFORMANCE EVALUATION

In order to evaluate the performance of DeepCog in realistic settings, we consider the mobile network infrastructure of a major operator in a large metropolitan region. The area under study covers around 100 km² with a resident population of more than 2 millions, and is surrounded by a conurbation of 11 millions inhabitants who often commute to it. We run DeepCog on real-world measurement data of an operator with a market share of 35% in the target region that captures the traffic generated by millions of users. The data were collected by monitoring the GPRS Tunneling Protocol (GTP) via dedicated probes deployed at the network gateway and the classification of IP flows into services was performed via Deep Packet Inspection (DPI) with proprietary models developed by the network operator. The traffic demands, expressed in bytes, refer to individual mobile services; they are aggregated at the antenna sector level and over intervals of 5 minutes. The demands capture the highly heterogeneous and time-varying loads that characterize real-world mobile network deployments, with differences in the offered traffic volume of up to two orders of magnitude between antenna sectors.

We consider services that belong to different categories, including video streaming, messaging and social networks.

These services impose a broad range of requirements into the network. For instance, video streaming services are consumed ubiquitously, they have significant bandwidth requirements (*e.g.*, up to 6 Mbps for a 1080p video in YouTube [36]) as well as latency constraints to avoid interruptions while playing the video [37], and they also require significant data center resources from the server side [36]. Messaging services, instead, have a large component of uplink traffic; overall, it has a rather relaxed requirements in terms of bandwidth and latency, as it does not involve any interactive communications. Social networks need considerable bandwidth and with some latency requirements; according to [38], a bandwidth of 8 Mbps and an access delay not exceeding 100 ms are required to achieve a high service quality. Beyond the specific requirements of each service, what really matters in the context of this paper is the fact that these are services of a very different nature with highly diverse requirements, and therefore they are likely to be served by different slices in a mobile network supporting network slicing.

We employ the measurement data to design three case studies combining several popular mobile services and different classes of network datacenters⁴. Each class is defined by the network location and number of served eNodeBs, ranging from centralized datacenters located in the core and serving many eNodeBs to more distributed ones located in the edge and serving a smaller number of eNodeBs. By selecting a diverse set of case studies, we can assess the DeepCog flexibility serving heterogeneous NFV scenarios, comprising different services and datacenter classes (C-RAN, MEC and core). In a first case study, we consider that a slice is instantiated for the incumbent video streaming service, *i.e.*, YouTube, at C-RAN datacenters in the target metropolitan area, each located in proximity of the radio access and performing baseband processing and scheduling for around ten eNodeBs. In the second case study, we look into Mobile Edge Computing (MEC) datacenters that handle the traffic of around 70 eNodeBs each, where a dedicated slice accommodates the traffic generated by Snapchat, a favored messaging app. The third case study focuses on a network slice dedicated to social network services provided by Facebook that are run at a core network datacenter controlling all 470 4G eNodeBs in the target metropolitan area.

The three case studies cover applications with diverse requirements in terms of bandwidth and latency; also, they entail very different spatiotemporal dynamics of the mobile traffic, as the considered services feature different loads and activity peaks [40]. In addition, the datacenter classes we consider have dissimilar geographical coverage and aggregated traffic volumes, as they serve the demands associated to a variable number of antennas, from ten to several hundreds. Overall, the three case studies considered for the DeepCog's performance evaluation are very useful to understand the effect of network slicing on the network operation costs. In fact, our results illustrate for the first time the impact of the slice isolation requirements –critical to future softwarized networks– on services that have a dominant role in today's traffic and are

³Fig. 4 shows results for α in the range [0.5, 3], and two exemplary T_h values, 5 and 30 minutes. Similar curves characterize all α values and prediction horizons (up to 8 hours) we tested.

⁴The internal organization of the mobile network – hence the demand recorded at each datacenter – is inferred by adopting the methodology proposed in [39].

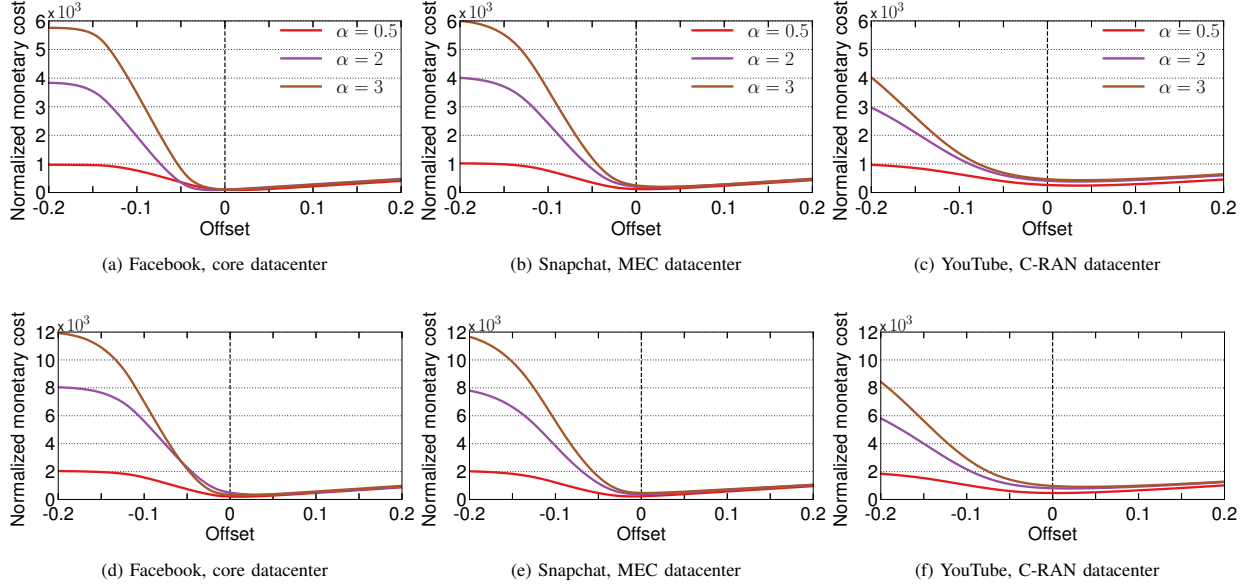


Fig. 4. Monetary cost (aggregated over time and normalized by the cost of one capacity unit) incurred when the overprovisioning level is shifted from that selected by DeepCog (at the abscissa origin). Each plot refers to one case study, *i.e.*, a combination of (i) mobile service associated to a dedicated slice and (ii) datacenter type. Top row: $T_h = 5$ minutes, bottom row: $T_h = 30$ minutes.

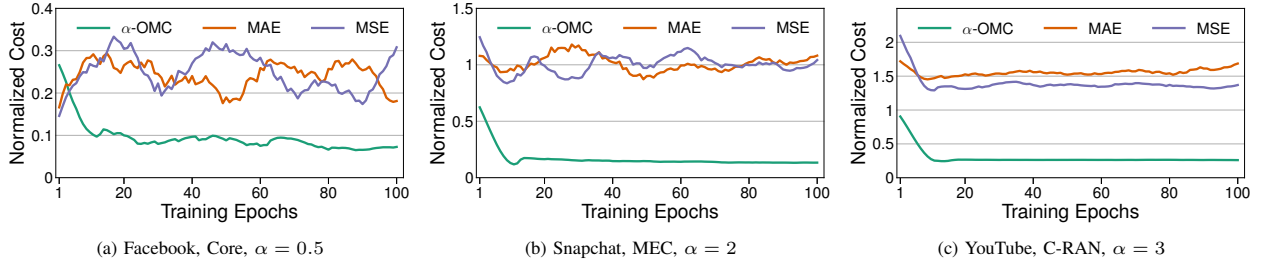


Fig. 5. Average cost versus the learning epochs, when the DeepCog neural network architecture is trained with α -OMC, MSE and MAE loss functions.

expected to keep playing a very relevant role in future mobile networks. In our tests, we do not parametrize different excess resource costs for each datacenter nor different SLA violation penalties for diverse services, leading to homogeneous α settings (*i.e.*, $\alpha_{sj} = \alpha$ for all $s \in \mathcal{S}, j \in \mathcal{M}$). However, since DeepCog provisions resources for each slice and each datacenter independently, by evaluating different α values our results provide insights on the behavior of heterogeneous α settings as well.

As discussed in Section III, DeepCog outputs a capacity forecast within a variable time-horizon T_h . We measure this time in the number of steps it comprises, where each step corresponds to the 5 mins granularity of our measurement data. In our evaluation T_h ranges from 5 minutes (which maps to a next-step prediction) to 8 hours (which corresponds to a forecast with a 96 time steps look-ahead). These are reasonable values in our context, since resource reallocation updates in the order of minutes are typical for computational and memory resources in architectures implementing NFV [41], and are in line with those supported by any state-of-the-art Virtual Infrastructure Manager (VIM) [42]. Conversely, larger intervals are more suitable for operations involving manual intervention, *e.g.*, spectrum leasing.

In all cases, we use the previous 30 minutes of traffic (*i.e.*, $T_p = 6$) as the DeepCog input, arranged in a 47×10 matrix as an input. This configuration proved to yield the best results when confronted to a number of other design strategies for the input that we explored, including longer, shorter, or non-continuous historical data time intervals. Capacity is predicted in terms of bytes of traffic, which is a reasonable metric to capture for resource utilization in actual virtual network functions [43], and is independent of the exact type of resources relevant for the mobile operator in each case study. We employ two months of mobile traffic data for training, two weeks of data for validation and another two for the actual experiments. This setting is also used for all benchmark approaches. All results are derived with a high level of confidence and low standard deviation.

A. Gain over state-of-the-art traffic predictors

We first focus on the particular case of next-step prediction, *i.e.*, $T_h = 5$ minutes, as this benchmark lets us compare our framework against state-of-the-art solutions that can only perform a forecast for the following time interval. We compare DeepCog against four benchmarks: (i) a naive technique that forecasts the future offered load by replicating the demand

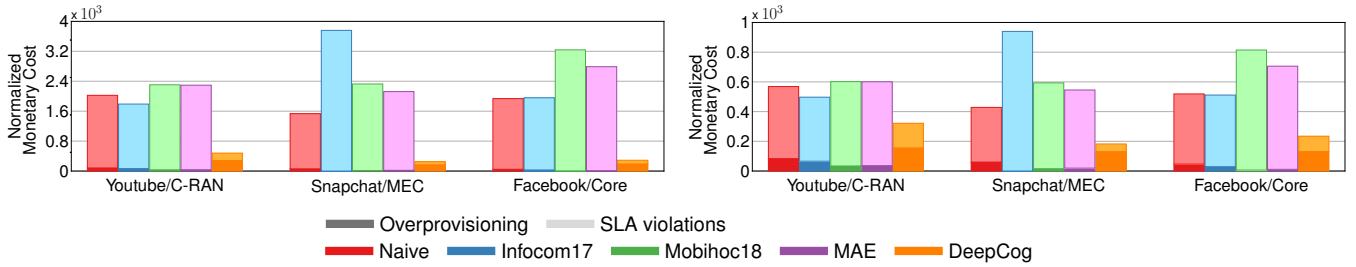


Fig. 6. Comparative evaluation of DeepCog with four benchmarks in three representative case studies. The monetary cost (normalized by the cost of one capacity unit) incurred by the operator is split into costs due to overprovisioning (dark) and SLA violations (light). Left: $\alpha = 2$. Right: $\alpha = 0.5$.

recorded at the same time during the previous week; (ii) the first approach proposed to predict mobile traffic based on a deep learning structure, referred to as *Infocom17* [8]; (iii) a recent solution for mobile network demand prediction that leverages a more complex deep neural network, referred to as *MobiHoc18* [21]; (iv) a reduced version of DeepCog, which replaces α -OMC with a legacy Mean Absolute Error (MAE) loss function⁵.

The results achieved in our three reference case studies by DeepCog and by the four benchmarks above are shown in Fig. 6. The plots report the normalized monetary cost for the operator, broken down into the expenses for unnecessary resource allocation (*i.e.*, overprovisioning) and fees for unserved demands (*i.e.*, SLA violations). We observe that DeepCog yields substantially lower costs than all other solutions. Indeed, the cost incurred by DeepCog for $\alpha = 2$ ranges between 15% (Facebook/Core) and 27% (Youtube/C-RAN) of the cost provided by the *best* competitor, depending on the case study. *Infocom17*, as all other benchmarks, targets mobile network traffic prediction, whereas DeepCog aims at forecasting capacity. As a result, DeepCog balances overprovisioning and SLA violations so as to minimize operation expenses, while *Infocom17* is oblivious to such practical resource management considerations. In other words, legacy predictors follow as closely as possible the general trend of the time series and allocate resources based on their prediction, which leads to systematic SLA violations that are not acceptable from a market viewpoint and determine huge fees for the operator. Instead, DeepCog selects the appropriate level of overprovisioning that, by suitably overestimating the offered load, minimizes monetary penalties (see Fig. 4). Indeed, even when choosing a low value such as $\alpha = 0.5$, which inflicts a small penalty for a SLA violation, the cost incurred by DeepCog is 64% of that incurred by the best performing benchmark.

B. Comparison with overprovisioned traffic prediction

In the light of the above results, a more reasonable approach to resource allocation could be to consider a traditional mobile traffic prediction as a basis, and adding some *overprovisioning offset* on top of it. In order to explore the effectiveness of such an approach, we design and implement several variants to MAE, as follows.

⁵We also experimented with other popular loss functions, *e.g.*, Mean Squared Error (MSE), with comparable results, omitted for space reasons.

A first variant adds an *a-posteriori* constant overprovisioning offset to the MAE output. This strategy, referred to as *MAE-post*, requires selecting a value of the static offset, which is then added to the predicted traffic. We dimension the offset as a certain percentage of the peak traffic activity observed in the whole historical data, and set it at 5%, which we deem a reasonable value in presence of a decently accurate prediction. Alternatively, we also consider a best-case version of this solution, named *MAE-post-best*, where an a-posteriori overprovisioning is chosen by performing an exhaustive search over all possible offset values and selecting the one that minimizes the loss function $\ell(\cdot)$.

A second variant accounts for some level of overprovisioning in a *preemptive* fashion, by introducing the offset during the deep neural network training. To this end, the *MAE-pre* solution replaces the MAE loss function with a new loss function $\mathcal{O} + \frac{1}{M} \sum_{j \in \mathcal{M}} |c_s^j(t) - d_s^j(t)|$, where \mathcal{O} denotes the a-priori overprovisioning offset. Also in this case, we set \mathcal{O} equal to 5% of the peak traffic in the historical data. To compare against the best possible operation of this scheme, we also consider a *MAE-pre-best* variant where \mathcal{O} is set equal to the average overprovisioning level provided by DeepCog for the test period.

We remark that the *MAE-post-best* and *MAE-pre-best* approaches are oracles and not feasible in practice, since they require knowledge of the future to determine the best a-posteriori values for the offset and the value of \mathcal{O} , respectively. Yet, they provide a benchmark for comparing the performance of DeepCog against optimal solutions that rely on traditional mobile network traffic prediction.

Fig. 7 shows the relative performance of the four variants above with respect to that attained by DeepCog, for $T_h = 5$ min (left) and $T_h = 45$ min (right). The figure shows the overprovisioned capacity, unserved traffic, and total economic cost incurred by the operator relative to the performance offered by DeepCog (in percentage). For $T_h = 5$ min, the results highlight how using a static overprovisioning in combination with a traditional traffic prediction is largely suboptimal, both when the additional offset is considered preemptively or a-posteriori. Indeed, the two practical solutions considered, *i.e.*, *MAE-post* and *MAE-pre*, cause SLA violations that are two- to three-fold more frequent than that incurred into by DeepCog, resulting in an economic cost that is 140% to 400% higher. Interestingly, even when parametrized with the best possible offsets, the approaches based on legacy traffic prediction

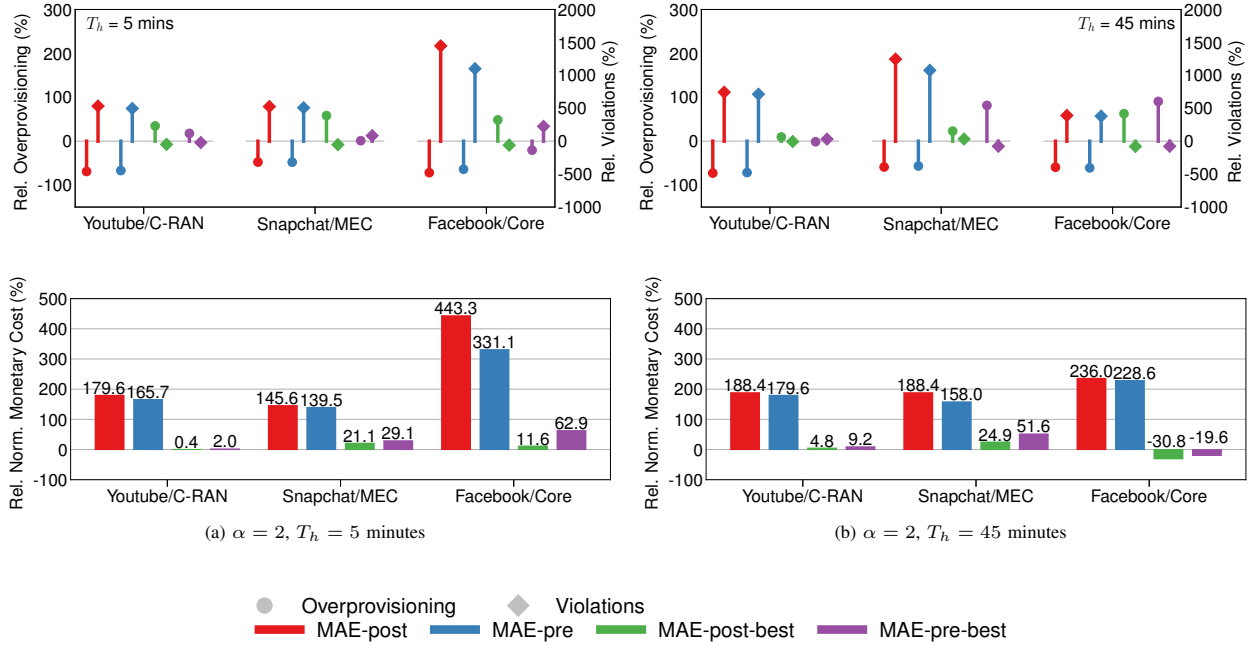


Fig. 7. Relative performance of overprovisioned traffic predictors, expressed as a percent of the cost attained by DeepCog. Top: relative overprovisioning and SLA violations. Bottom: relative monetary cost. Results refer to $\alpha = 2$ and prediction horizons of 5 (left) and 45 (right) minutes.

cannot match the performance of DeepCog: MAE-post-best and MAE-pre-best dramatically reduce the penalties of their viable counterparts, yet lead to monetary costs that are up to 60% higher than those of DeepCog.

The results for $T_h = 45$ min⁶, show that the above considerations hold across different values of the prediction horizon. The advantage over feasible overprovisioned traffic predictors such as MAE-pre or MAE-post is aligned with that observed under a next-step prediction, as such solutions increase the overall cost by 188% to 236%. When considering a long horizon of 45 minutes, oracle methods based on overprovisioning like (*i.e.*, MAE-pre-best and MAE-post-best) can outperform DeepCog, further reducing the operator cost by 19% to 30%. This is due to the fact that, when prediction must be performed with significant time advance, the accuracy of DeepCog cannot be as high as an oracle that knows future demand and has hence a significant advantage. However, even under such conditions DeepCog performs almost as good as the oracles or better.

We conclude that traffic predictors – no matter how they are enhanced – are not appropriate for the capacity forecast problem, for the simple reason that they are designed for a different purpose. Indeed, they ignore the economic penalties incurred by SLA violations, and this limits drastically their ability to address this problem. Strategies that rely on integrating such costs into the solution after the traffic prediction is performed are largely suboptimal.

⁶Note that, in order to perform a fair comparison, we had to extend the MAE policy to compute the average absolute error on each time slot in the $[t, \dots, t + T_h]$ interval.

C. Controlling resource allocation trade-offs with the α parameter of DeepCog

As discussed in Section IV, DeepCog addresses a fundamental trade-off between overprovisioning and SLA violations, aiming to find the best possible compromise between the two. An operator is given the flexibility of choosing the desired operation point within this trade-off, by suitably setting the α parameter. In the following, we carry out an extensive analysis of the trade-off between overprovisioning of resources and failing to meet service demands. This study is conducted for a large number of practical scenarios that extend the original three case studies considered in the comparative analysis. Specifically, we select five different network slices, dedicated to the same number of popular mobile services: the three we already studied, *i.e.*, YouTube, Facebook, and Snapchat, plus iTunes and Instagram. We then investigate the performance of DeepCog when such slices are deployed at the three classes of datacenters introduced before, *i.e.*, at the C-RAN, MEC, and network core. Overall, this leads to 15 distinct scenarios.

Fig. 8a shows results in all of the above settings under different economic strategies that are reflected by the α parameter of the loss function $\ell(\cdot)$. Configurations range from policies that prioritize minimizing overprovisioning over avoiding SLA violations ($\alpha = 0.5$) to others that strictly enforce the SLAs at the price of allocating additional resources ($\alpha = 5$). The plots tell apart the contribution of the two components that contribute to the total monetary cost: overprovisioning is expressed as a percentage of the actual demand, and SLA violations are measured as a percentage of the time slots in the test period. As expected, higher α values reduce the number SLA violations, as they become increasingly expensive; this occurs at the cost of provisioning additional capacity, which

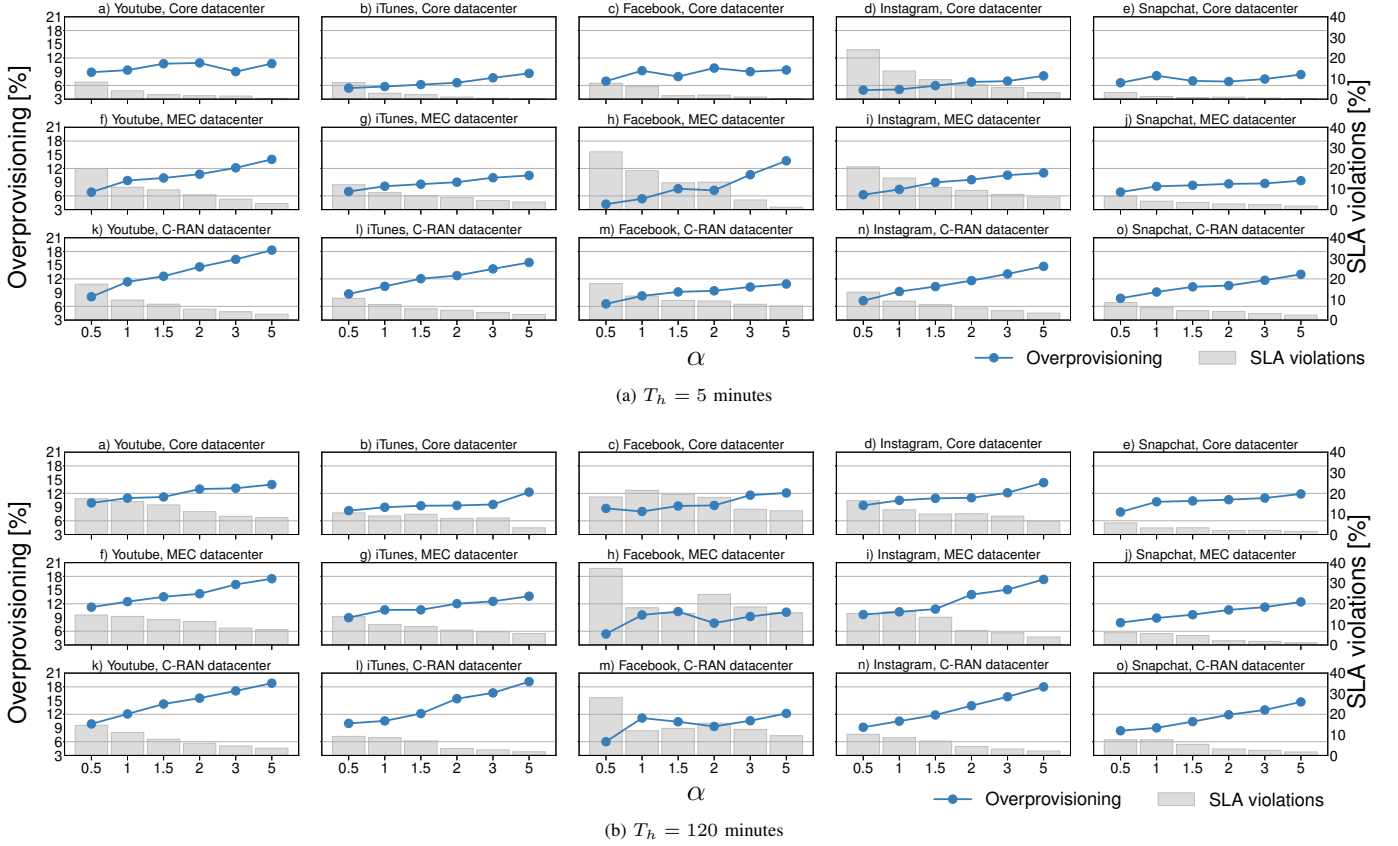


Fig. 8. Tradeoff between resource overprovisioning (expressed as a percentage of the actual demand) and SLA violation (expressed as a percentage of time slots), as a function of the α parameter. Results refer to 15 different scenarios, and two values of the prediction horizon T_h , i.e., 5 minutes (a) and 120 minutes (b).

becomes instead cheaper in proportion⁷. The trend is consistent across all scenarios, confirming that α effectively drives resource orchestration towards the desired operation point.

Our analysis also reveals that the level of overprovisioning grows in most cases as one moves from datacenters in the network core outwards. This trend applies across slices, and is due to the fact that more centralized datacenters serve an increasingly aggregate traffic that is generally less noisy and easier to predict. Under such conditions, DeepCog needs a reduced additional capacity to limit unserved demands: as a result, SLA violations are often lower at core datacenters.

Fig. 8a refers to a short-term prediction for $T_h = 5$ minutes, however the same trends discussed above are confirmed for larger prediction horizons. For instance, Fig. 8b reports the same results for $T_h = 120$ minutes. The only remarkable difference is that overprovisioning and SLA violations are higher than in the case of a 5-minute prediction, as forecasting on larger time horizons is obviously harder. Yet, the impact of α is equivalent to that observed for $T_h = 5$ minutes. We analyze in more detail DeepCog's performance as a function of T_h in the next section.

Overall, the results presented above show that DeepCog finds good trade-offs between resource overprovisioning and

SLA violations in very different cost settings across slices and datacenter types. Since each DeepCog instance for a slice at a datacenter runs independently, this shows that DeepCog will grant good performance also in scenarios where resource costs may differ across datacenters, e.g., due to diverse operation and management costs in urban and rural facilities.

D. Long-term capacity prediction with DeepCog

DeepCog aims at forecasting the (constant) capacity that should be allocated over a long-term horizon, so as to minimize the monetary cost incurred by the operator. As discussed in Section III, this is particularly useful in practical settings where the NFV technology imposes limits on the frequency upon which resources can be reallocated. In this section, we thoroughly study how the performance of DeepCog varies with the prediction horizon.

Fig. 9 summarizes the overall trend of the monetary cost incurred by DeepCog, as the periodicity of the reconfiguration opportunities ranges from 5 minutes to 8 hours. The plots outline a diversity of scenarios, combining different datacenter classes (C-RAN, MEC, and core) and relative expenses of overprovisioning and SLA violations (α equal to 0.5, 2, and 5). The results correspond to the case where one slice is dedicated to the traffic generated by YouTube, but equivalent behaviors were observed for the other services. In all settings, the cost grows with the prediction horizon, which, as already mentioned, is largely expected. What is less expected, however, is

⁷Sporadic disruptions in the monotonicity of the cost curves are due to the inherent randomness of the measurement data; indeed, the data correspond to a specific time period and it may show some biases that would not be observed over different (or longer) time periods.

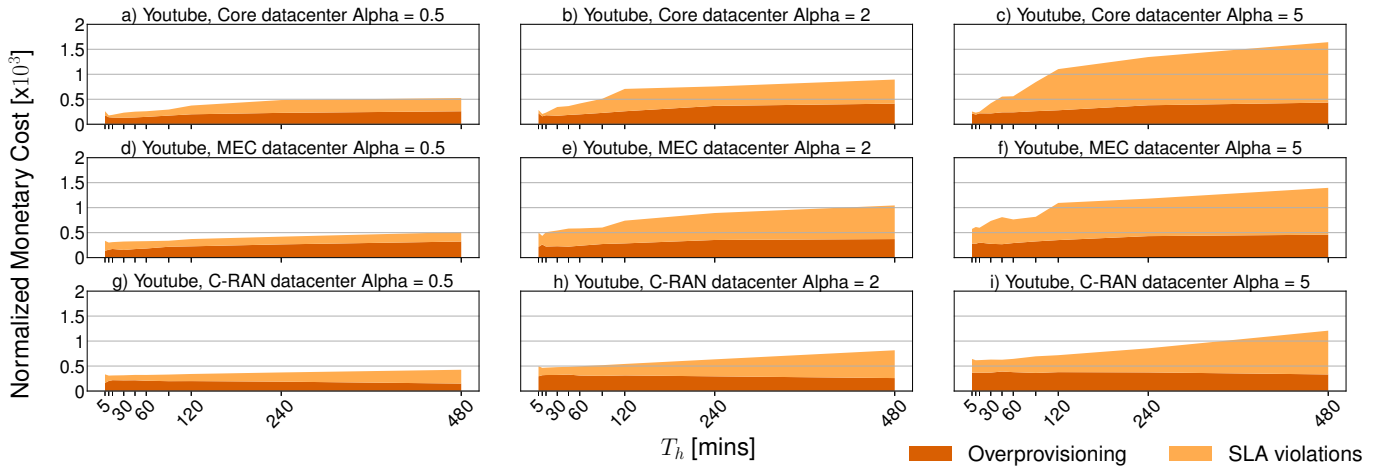


Fig. 9. Monetary cost (normalized by the cost of one capacity unit) incurred by the operator, versus the prediction time horizon T_h . The plots refer to different combinations of datacenter class and economic strategies modeled by α , for a slice dedicated to the YouTube mobile service.

the quasi-linear relationship between the cost and T_h . This is a very important result, as it shows that even if we increase the intervals for resource reallocation (*i.e.*, the time horizon), the economic expenses of the operator remain bounded and do not skyrocket (as they would if the growth was, *e.g.*, exponential). The result thus demonstrates the efficiency of DeepCog in limiting the unavoidable increased penalty associated to forecasting long-term capacity: as an indicative figure, the cost is roughly increased by two when moving from a 5-minute prediction to one that spans the following 8 hours which is a very reasonable factor.

The impact of the other system parameters is in line with our previous analysis: higher monetary fees for SLA violations (*i.e.*, higher α values) lead to increased costs, whereas the performance is comparable across resource allocations over different classes of datacenter (C-RAN, MEC and core), each corresponding to different traffic volumes. It is nonetheless interesting to note that the property of a linear growth of the cost over T_h is preserved under any combination of such parameters.

Fig. 9 also offers a breakdown of the overall monetary costs into the two contributions (overprovisioning and SLA violations). Violations of SLAs yield substantially higher absolute costs and dominate the increase of total cost with T_h ; the effect is clearly stronger for higher values of α . A more detailed view that highlights the exact evolution of the two cost components as a function of T_h is provided in Fig. 10, showing that both contribute to increasing costs over longer-term forecasts. However, and interestingly, the dynamics of the two components with T_h are diverse depending on the system settings such as the datacenter class and the value of α . The common trend here is that the penalty associated with both overprovisioning and SLA violations is fairly stable when the horizon is increased from 5 minutes up to two hours. For forecasts beyond two hours, however, these fees (one of the two or both) tend to increase substantially with T_h .

We ascribe these behaviors to (i) the relationship between T_h and the timescale of temporal fluctuations in the input demand, and (ii) the way DeepCog reacts to the problem

of devising a capacity forecast, which becomes harder for larger T_h values. The first point relates to the characteristics of the input data (see Fig. 11 for illustrative examples of the temporal oscillation of the service demand). For very large T_h (above 120 minutes) the prediction task performed by DeepCog resorts to an “envelope” of the demand that accommodates the peak over the T_h period. This means that for those times where demand is below the peak, we incur a high level of overprovisioning that increases the resulting cost. In contrast, smaller T_h allow to adapt the capacity forecasting to the actual demand at each point in time, providing an advantage in terms of cost. The second point relates to the behavior and performance of DeepCog under large T_h values. DeepCog aims at providing a similar level of overprovisioning over time, as exemplified by the top three plots of Fig. 11. For large T_h values this yields increased SLA violations, since the oscillations make it more likely that the constant capacity falls below the demand curve at some point during T_h . Additionally, larger T_h values make the prediction task inherently harder, which further contributes to increasing the SLA violations costs.

VI. CONCLUSIONS

In this paper we presented and evaluated DeepCog, an original data analytics tool for the cognitive management of resources in sliced 5G networks. DeepCog tackles the novel problem of capacity forecasting, whose solution is key to the sustainable operation of future multi-tenant mobile networks. Inspired by recent advances in deep learning for image and video processing, DeepCog hinges upon a deep neural network structure, which analyzes antenna-level demand snapshots for different services in order to provide a prediction of the resources that the operator has to allocate to accommodate the future load. The operation is performed for individual mobile services separately, and over a configurable time horizon. At the core of DeepCog there is α -OMC, a new and customized loss function that drives the deep neural network training so as to minimize the monetary cost contributed by two main deployment fees, *i.e.*, overprovisioning and SLA violation. Ours

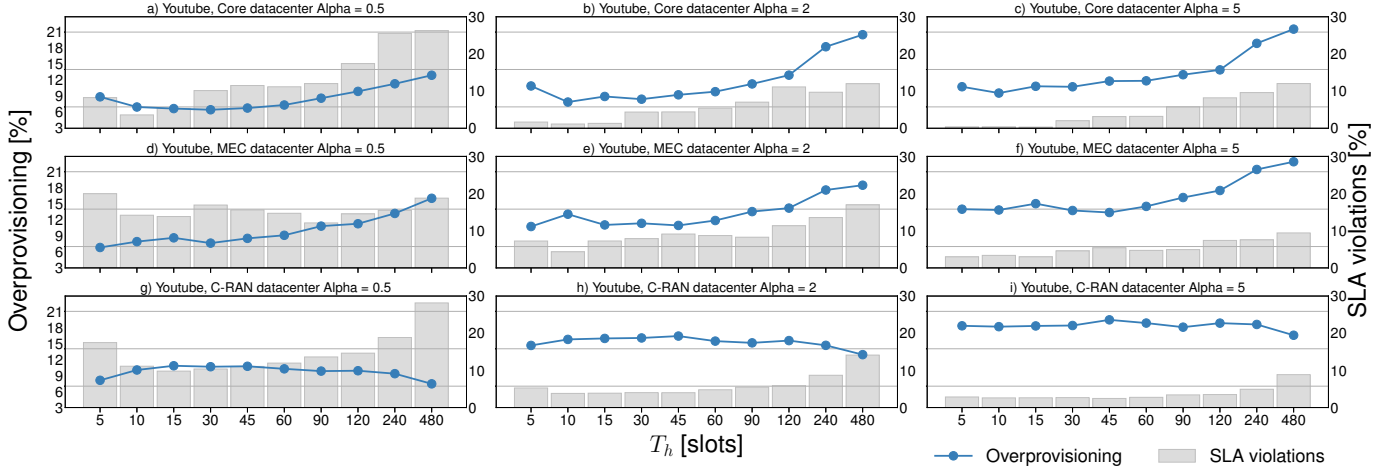


Fig. 10. Breakdown of monetary costs into two contributions: (i) overprovisioning (expressed as a percentage of the actual demand) and (ii) SLA violations (expressed as a percentage of time slots), in the scenarios of Fig. 9.

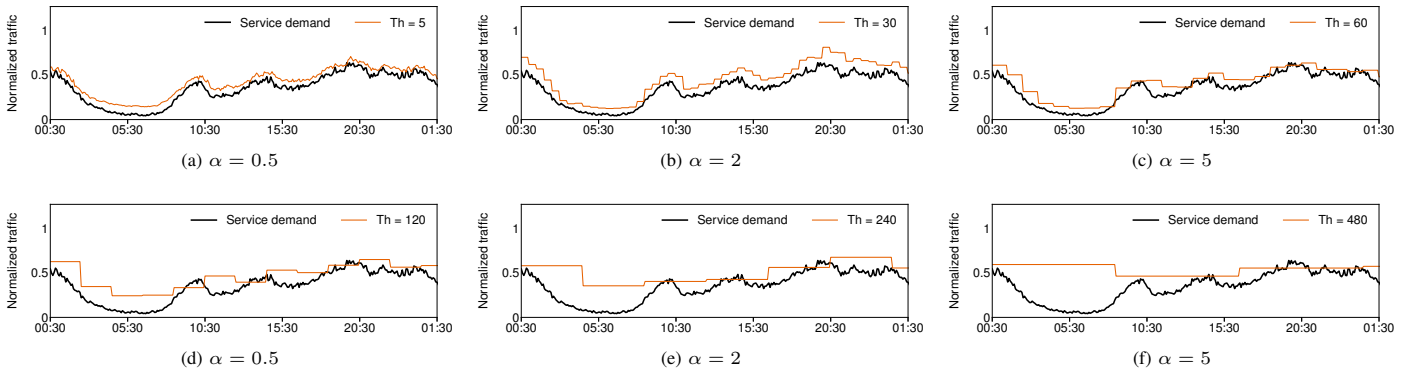


Fig. 11. Illustrative examples of the capacity forecast returned by DeepCog behavior under different prediction time horizons. The scenario refers to a network slice dedicated to the YouTube mobile service that is deployed at a core datacenter, under $\alpha = 2$.

is, to the best of our knowledge, the only work to date where a deep learning architecture is explicitly tailored to the problem of anticipatory resource orchestration in mobile networks. The solution presented in this paper thus represents a first attempt to integrate data analytics based on machine learning into an overall cognitive management framework. Thorough empirical evaluations with real-world metropolitan-scale data show the substantial advantages granted by DeepCog over state-of-the-art predictors and other automated orchestration strategies, providing a first analysis of the practical costs of heterogeneous network slice management across a variety of case studies.

ACKNOWLEDGMENTS

The work of University Carlos III of Madrid was supported by H2020 5G-TOURS project (grant agreement no. 856950). The work of NEC Laboratories Europe was supported by H2020 5G-TRANSFORMER project (grant agreement no. 761536) and 5GROWTH project (grant agreement no. 856709). The work of CNR-IEIT was partially supported by the ANR CANCAN project (ANR-18-CE25-0011).

REFERENCES

- [1] P. Rost *et al.*, “Network Slicing to Enable Scalability and Flexibility in 5g Mobile Networks,” *IEEE Communications Magazine*, vol. 55, no. 5, pp. 72–79, May 2017.
- [2] X. Foukas *et al.*, “Network Slicing in 5G: Survey and Challenges,” *IEEE Communications Magazine*, vol. 55, no. 5, pp. 94–100, May 2017.
- [3] A. de la Oliva *et al.*, “5G-TRANSFORMER: Slicing and Orchestrating Transport Networks for Industry Verticals,” *IEEE Communications Magazine*, vol. 56, no. 8, pp. 78–84, Aug. 2018.
- [4] 5G-PPP, “The 5G Infrastructure Association. Pre-structuring Model, version 2.0,” Nov. 2017.
- [5] ETSI, NFVGS, “Network Function Virtualization (NFV) Management and Orchestration,” *NFV-MAN*, vol. 1, Dec. 2014.
- [6] ETSI, “Open Source MANO (OSM) Project.” [Online]. Available: <https://osm.etsi.org/>
- [7] The Linux Foundation, “ONAP, Open Network Automation Framework.” [Online]. Available: <https://www.onap.org>
- [8] J. Wang *et al.*, “Spatiotemporal modeling and prediction in cellular networks: A big data enabled deep learning approach,” in *Proc. of IEEE INFOCOM*, Atlanta, GA, USA, May 2017, pp. 1–9.
- [9] C. Zhang *et al.*, “Deep Learning in Mobile and Wireless Networking: A Survey,” *arXiv:1803.04311 [cs.NI]*, Mar. 2018.
- [10] M. Wang *et al.*, “Machine Learning for Networking: Workflow, Advances and Opportunities,” *IEEE Network*, vol. 32, no. 2, pp. 92–99, Mar. 2018.
- [11] J. X. Salvat *et al.*, “Overbooking Network Slices Through Yield-driven End-to-end Orchestration,” in *Proc. of ACM CoNEXT*, Heraklion, Greece, Dec. 2018, pp. 353–365.

- [12] J. J. Ayala *et al.*, "vRAIn: A Deep Learning Approach Tailoring Computing and Radio Resources in Virtualized RANs," in *Proc. of ACM MobiCom*, Los Cabos, Mexico, Nov. 2019, pp. 1–16.
- [13] C. Gutterman *et al.*, "RAN resource usage prediction for a 5G slice broker," in *Proc. of ACM Mobihoc*, Catania, Italy, Jul. 2019.
- [14] M. Joshi and T. H. Hadi, "A Review of Network Traffic Analysis and Prediction Techniques," *arXiv:1507.05722 [cs.NI]*, Jul. 2015.
- [15] F. Xu *et al.*, "Big Data Driven Mobile Traffic Understanding and Forecasting: A Time Series Approach," *IEEE Transactions on Services Computing*, vol. 9, no. 5, pp. 796–805, Sep. 2016.
- [16] M. Zhang *et al.*, "Understanding Urban Dynamics From Massive Mobile Traffic Data," *IEEE Transactions on Big Data*, vol. 5, no. 2, pp. 266–278, Nov. 2017.
- [17] S. T. Au *et al.*, "Automatic forecasting of double seasonal time series with applications on mobility network traffic prediction," *JSM Proceedings, Business and Economic Statistics Section*, Jul. 2011.
- [18] R. Li *et al.*, "The prediction analysis of cellular radio access network traffic: From entropy theory to networking practice," *IEEE Communications Magazine*, vol. 52, no. 6, pp. 234–240, Jun. 2014.
- [19] M. Z. Shafiq *et al.*, "Characterizing and modeling internet traffic dynamics of cellular devices," in *Proc. of ACM SIGMETRICS*, San Jose, CA, USA, Jun. 2011, pp. 265–276.
- [20] A. Y. Nikraves *et al.*, "An Experimental Investigation of Mobile Network Traffic Prediction Accuracy," *Services Transactions on Big Data*, vol. 3, no. 1, pp. 1–16, Jan. 2016.
- [21] C. Zhang *et al.*, "Long-Term Mobile Traffic Forecasting Using Deep Spatio-Temporal Neural Networks," in *Proc. of ACM MobiHoc*, Los Angeles, CA, USA, Jun. 2018, pp. 231–240.
- [22] S. Ntalampiras *et al.*, "Forecasting mobile service demands for anticipatory MEC," in *Proc. of IEEE WoWMoM*, Chania, Greece, Jun. 2018, pp. 14–19.
- [23] D. Bega *et al.*, "DeepCog: Cognitive Network Management in Sliced 5G Networks with Deep Learning," in *Proc. of IEEE INFOCOM*, Paris, France, May 2019, pp. 1–9.
- [24] —, "Alfa-OMC: cost-aware deep learning for mobile network resource orchestration," in *Proc. of IEEE INFOCOM NI Workshop*, Paris, France, May 2019, pp. 1–9.
- [25] K. Simonyan *et al.*, "Very deep convolutional networks for large-scale image recognition," *arXiv:1409.1556*, Sep. 2014.
- [26] Y. LeCun, "Generalization and network design strategies," *Connectionism in perspective*, vol. 19, pp. 143–155, Jun. 1989.
- [27] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. of IEEE CVPR*, Boston, MA, USA, Jun. 2015, pp. 1–9.
- [28] G. E. Dahl *et al.*, "Improving deep neural networks for LVCSR using rectified linear units and dropout," in *Proc. of IEEE ICASSP*, Vancouver, Canada, May 2013, pp. 8609–8613.
- [29] M. W. Gardner *et al.*, "Artificial neural networks (the multilayer perceptron) - A review of applications in the atmospheric sciences," *Atmospheric environment*, vol. 32, no. 14–15, pp. 2627–2636, Aug. 1998.
- [30] F. Chollet, *Deep learning with Python*. Manning Publications Co, 2018.
- [31] D. P. Kingma *et al.*, "Adam: A method for stochastic optimization," *arXiv:1412.6980*, Dec. 2014.
- [32] A. Furno *et al.*, "A Tale of Ten Cities: Characterizing Signatures of Mobile Traffic in Urban Areas," *IEEE Transactions on Mobile Computing*, vol. 16, no. 10, pp. 2682–2696, Oct. 2017.
- [33] J. Paparrizos *et al.*, "k-Shape: Efficient and Accurate Clustering of Time Series," in *Proc. of ACM SIGMOD*, Melbourne, Victoria, Australia, Jun. 2015, pp. 1855–1870.
- [34] I. Borg *et al.*, "Modern Multidimensional Scaling: Theory and Applications," *Journal of Educational Measurement*, vol. 40, no. 3, pp. 277–280, Sep. 2003.
- [35] H. W. Kuhn, "The Hungarian method for the assignment problem," *Naval Research Logistics Quarterly*, vol. 2, no. 1–2, pp. 83–97, Mar. 1955.
- [36] D. K. Krishnappa *et al.*, "DASHing YouTube: An analysis of using DASH in YouTube video service," in *Proc. of IEEE LCN*, Sydney, NSW, Australia, 2013, pp. 407–415.
- [37] G. Dimopoulos *et al.*, "Measuring video QoE from encrypted traffic," in *Proc. of ACM IMC*, Santa Monica, CA, USA, Nov. 2016, pp. 513–526.
- [38] P. Casas *et al.*, "Qomson-on the analysis of traffic and quality of experience in mobile online social networks," in *Proc. of IEEE EuCNC*, Paris, France, 2015, pp. 471–475.
- [39] C. Marquez *et al.*, "How Should I Slice My Network?: A Multi-Service Empirical Evaluation of Resource Sharing Efficiency," in *Proc. of ACM MobiCom*, New Delhi, India, Nov. 2018, pp. 191–206.
- [40] —, "Not All Apps Are Created Equal: Analysis of Spatiotemporal Heterogeneity in Nationwide Mobile Service Usage," in *Proc. of ACM CoNEXT*, Nov. 2017, pp. 180–186.
- [41] V. Sciancalepore *et al.*, "Mobile traffic forecasting for maximizing 5G network slicing resource utilization," in *Proc. of IEEE INFOCOM*, Atlanta, GA, USA, May 2017, pp. 1–9.
- [42] J. Gil Herrera *et al.*, "Resource Allocation in NFV: A Comprehensive Survey," *IEEE Transactions on Network and Service Management*, vol. 13, no. 3, pp. 518–532, Sep. 2016.
- [43] J.-J. Kuo *et al.*, "Service chain embedding with maximum flow in software defined network and application to the next-generation cellular network architecture," in *Proc. of IEEE INFOCOM*, Atlanta, GA, USA, May 2017, pp. 1–9.



Dario Bega is a research assistant with IMDEA Networks Institute in Madrid and a Ph.D. student in Telematic Engineering at University Carlos III of Madrid (UC3M). He received his B.Sc. (2010) and M.Sc. (2013) degrees in Telecommunication Engineering from the University of Pisa, Italy. His research work focuses on wireless networks, multi-tenancy and machine learning approaches for 5G.



Marco Gramaglia is a post-doc researcher at University Carlos III of Madrid (UC3M), where he received M.Sc. (2009) and Ph.D. (2012) degrees in Telematics Engineering. He held post-doctoral research positions at ISMB (Italy), the CNR-IEIT (Italy) and IMDEA Networks (Spain). He was involved in EU projects and authored more than 40 papers appeared in international conference and journals.



Marco Fiore (S'05, M'09, SM'17) is a researcher at CNR-IEIT (Italy) a Royal Society visiting research fellow, and a Marie Curie fellow. He received a PhD degree from Politecnico di Torino (Italy) and a HDR degree from Université de Lyon (France). He was associate professor at INSA Lyon, (France), associate researcher at Inria (France), visiting researcher at Rice University (USA) and UPC, (Spain), and visiting research fellow at UCL (UK). His current research interests are in the fields of mobile networks, network traffic analytics and privacy.



Albert Banchs (M04-SM12) received his M.Sc. and Ph.D. degrees from the Polytechnic University of Catalonia (UPC-BarcelonaTech) in 1997 and 2002, respectively. He worked at ICSI (Berkeley, California) in 1997, Telefonica I+D (Madrid, Spain) in 1998 and NEC Laboratories (Heidelberg, Germany) from 1998 to 2003. Since 2003 he is with the University Carlos III of Madrid, Spain, where he is currently Full Professor and has a double affiliation as Deputy Director of the IMDEA Networks institute. Prof. Banchs has served in many conference TPCs and journal editorial boards, and is currently Editor for IEEE/ACM Transactions on Networking.



Xavier Costa-Pérez (M'06-SM'18) is Head of 5G Networks R&D and Deputy General Manager of the Security and Networking Research Division at NEC Laboratories Europe. His team contributes to products roadmap evolution as well as to European Commission projects and received several awards for successful technology transfers. Dr. Costa is a 5GPPP Technology Board member, has served on the Program Committees of several conferences (including IEEE Greencom, WCNC, and INFOCOM), published at top research venues and holds several patents. He received both his M.Sc. and Ph.D. degrees in Telecommunications from the Polytechnic University of Catalonia (UPC) in Barcelona and was the recipient of a national award for his Ph.D. thesis.